

ARTICLE

<https://doi.org/10.1038/s42003-019-0415-5>

OPEN

diffcyt: Differential discovery in high-dimensional cytometry via high-resolution clustering

Lukas M. Weber ^{1,2}, Malgorzata Nowicka^{1,2,3}, Charlotte Soneson ^{1,2,4} & Mark D. Robinson ^{1,2}

High-dimensional flow and mass cytometry allow cell types and states to be characterized in great detail by measuring expression levels of more than 40 targeted protein markers per cell at the single-cell level. However, data analysis can be difficult, due to the large size and dimensionality of datasets as well as limitations of existing computational methods. Here, we present *diffcyt*, a new computational framework for differential discovery analyses in high-dimensional cytometry data, based on a combination of high-resolution clustering and empirical Bayes moderated tests adapted from transcriptomics. Our approach provides improved statistical performance, including for rare cell populations, along with flexible experimental designs and fast runtimes in an open-source framework.

¹Institute of Molecular Life Sciences, University of Zurich, CH-8057 Zurich, Switzerland. ²SIB Swiss Institute of Bioinformatics, University of Zurich, CH-8057 Zurich, Switzerland. ³Present address: F. Hoffmann-La Roche AG, CH-4070 Basel, Switzerland. ⁴Present address: Friedrich Miescher Institute for Biomedical Research and SIB Swiss Institute of Bioinformatics, CH-4058 Basel, Switzerland. Correspondence and requests for materials should be addressed to M.D.R. (email: mark.robinson@imls.uzh.ch)

High-dimensional flow cytometry and mass cytometry (or CyTOF, for “cytometry by time-of-flight mass spectrometry”) characterize cell types and states by measuring expression levels of pre-defined sets of surface and intracellular proteins in individual cells, using antibodies tagged with either fluorochromes (flow cytometry) or heavy metal isotopes (mass cytometry). Modern flow cytometry systems allow simultaneous detection of more than 20 proteins per cell, in thousands of cells per second¹. In mass cytometry, the use of metal tags significantly reduces signal interference due to spectral overlap and auto-fluorescence, enabling detection of more than 40 proteins per cell in hundreds of cells per second^{1,2}. Recently, further increases in the number of detected proteins have been demonstrated using oligonucleotide-tagged antibodies and single-cell sequencing³; this has also been combined with single-cell RNA sequencing on the same cells^{4,5}.

The rapid increase in dimensionality has led to serious bottlenecks in data analysis. Traditional analysis by visual inspection of scatterplots (“manual gating”) is unreliable and inefficient in high-dimensional data, does not scale readily, and cannot easily reveal unknown cell populations¹. Significant efforts have been made to develop computationally guided or automated methods that do not suffer from these limitations. For example, unsupervised clustering algorithms are commonly used to define cell populations in one or more biological samples. Recent benchmarking studies have demonstrated that several clustering methods can accurately detect known cell populations in low-dimensional flow cytometry data⁶, and both major and rare known cell populations in high-dimensional data⁷. A further benchmarking study comparing supervised methods for inferring cell populations associated with a censored continuous clinical variable demonstrated good performance for two methods using data of moderate dimensionality⁸.

Several new methods have recently been developed for performing (partially) supervised analyses with the aim of inferring cell populations or states associated with an outcome variable in high-dimensional cytometry data, including *Citrus*⁹, *CellCnn*¹⁰, *cydar*¹¹, and a *classic* regression-based approach¹² (a similar regression-based approach was also recently described by ref. 13). However, these existing methods have a number of limitations. In particular: detected features from *Citrus* cannot be ranked by importance, and the ranking of detected cells from *CellCnn* cannot be interpreted in terms of statistical significance; rare cell populations are difficult to detect with *Citrus* and *cydar* (by contrast, *CellCnn* is optimized for analysis of rare populations); the response variable in the models for *Citrus* and *CellCnn* is the outcome variable, which makes it difficult to account for complex experimental designs; and *CellCnn* and *cydar* do not distinguish between “cell type” and “cell state” (e.g. functional) markers, which can make interpretation difficult.

Here, we present *diffcyt*, a new computational framework based on high-resolution unsupervised clustering together with supervised statistical analyses to detect cell populations or states associated with an outcome variable in high-dimensional cytometry data. The *diffcyt* methodology uses clustering to define cell populations, and empirical Bayes moderated tests adapted from transcriptomics for differential analysis. By default, our implementation uses the *FlowSOM* clustering algorithm¹⁴, given its strong performance and fast runtimes⁷. For the differential analyses, we use methods from *edgeR*^{15,16}, *limma*¹⁷, and *voom*¹⁸, which are widely used in the transcriptomics field; in addition, we include alternative methods adapted from the *classic* regression-based framework¹². In principle, other high-resolution clustering algorithms or differential testing methods could also be substituted. Our methods consolidate several aspects of functionality from existing methods. Similar to *cydar* and the *classic* regression

framework, our model specification uses the cytometry-measured features (cell population abundances or median expression of cell state markers within populations) as response variables, which enables analysis of complex experimental designs, including batch effects, paired designs, and continuous covariates. Linear contrasts enable testing of a wide range of hypotheses. Rare cell populations can easily be investigated, since the use of high-resolution clustering ensures that rare populations are unlikely to be merged into larger ones. In addition, as in *Citrus* and the *classic* regression framework, we optionally allow the user to split the set of protein markers into cell type and cell state markers. In this setup, cell type markers are used to define clusters representing cell populations, which are tested for differential abundance (DA); and median cell state marker signals per cluster are used to test for differential states (DS) within populations. We note that the underlying definitions of cell type and cell state can be challenging to apply to observed data, and may partially overlap. In general, cell type refers to relatively stable or permanent features of a cell’s identity, while cell state refers to transient features such as signaling or other functional states or the cell cycle^{19–21}. In our view, providing the ability to maintain this distinction within the methodology greatly improves biological interpretability, since the results can be directly linked back to known cell types or populations of interest¹². Finally, our methods have fast runtimes, enabling exploratory and interactive analyses.

Results

Overview and benchmarking strategy. Figure 1 provides a schematic overview of the *diffcyt* methodology (see Methods for further details), and Table 1 provides a summary of existing methods and their limitations. We demonstrate the performance of our methods using four benchmark datasets: two semi-simulated datasets (*AML-sim* and *BCR-XL-sim*) and two published experimental datasets (*Anti-PD-1* and *BCR-XL*). The semi-simulated datasets have been constructed by computationally introducing an artificial signal of interest (an in silico spike-in signal) into experimental data, thus reflecting the properties of real experimental data while also including a known ground truth that can be used to calculate statistical performance metrics. The experimental datasets, which do not contain a ground truth, are evaluated in qualitative terms. A complete description of all benchmark datasets is provided in Supplementary Note 1, and additional details on the comparisons with existing methods are included in Supplementary Note 2.

Improved performance for DA tests. The *AML-sim* dataset evaluates performance for detecting DA of rare cell populations (Figure 2). The dataset contains a spiked-in population of acute myeloid leukemia (AML) blast cells, in a comparison of 5 vs. 5 myeloid samples of otherwise healthy bone marrow mononuclear cells, which simulates the phenotype of minimal residual disease in AML patients (the data generation strategy is adapted from ref. 10, and uses original data from ref. 22). The simulation was repeated for two subtypes of AML (cytogenetically normal, CN; and core-binding factor translocation, CBF), and three thresholds of abundance for the spiked-in population (5%, 1%, and 0.1%). Figure 2a displays representative results for one subtype (CN) and one threshold (1%), for all *diffcyt* DA methods as well as *Citrus*, *CellCnn*, and *cydar* (complete results are included in Supplementary Fig. 1). Methods *diffcyt-DA-edgeR*, *diffcyt-DA-voom*, and *CellCnn* give the best performance; the *diffcyt* results can also be interpreted as adjusted *p*-values, enabling a standard statistical framework where a list of significant detected clusters is determined by specifying a cutoff for the false discovery rate (FDR). *diffcyt-DA-GLMM* has inferior error control at the given FDR

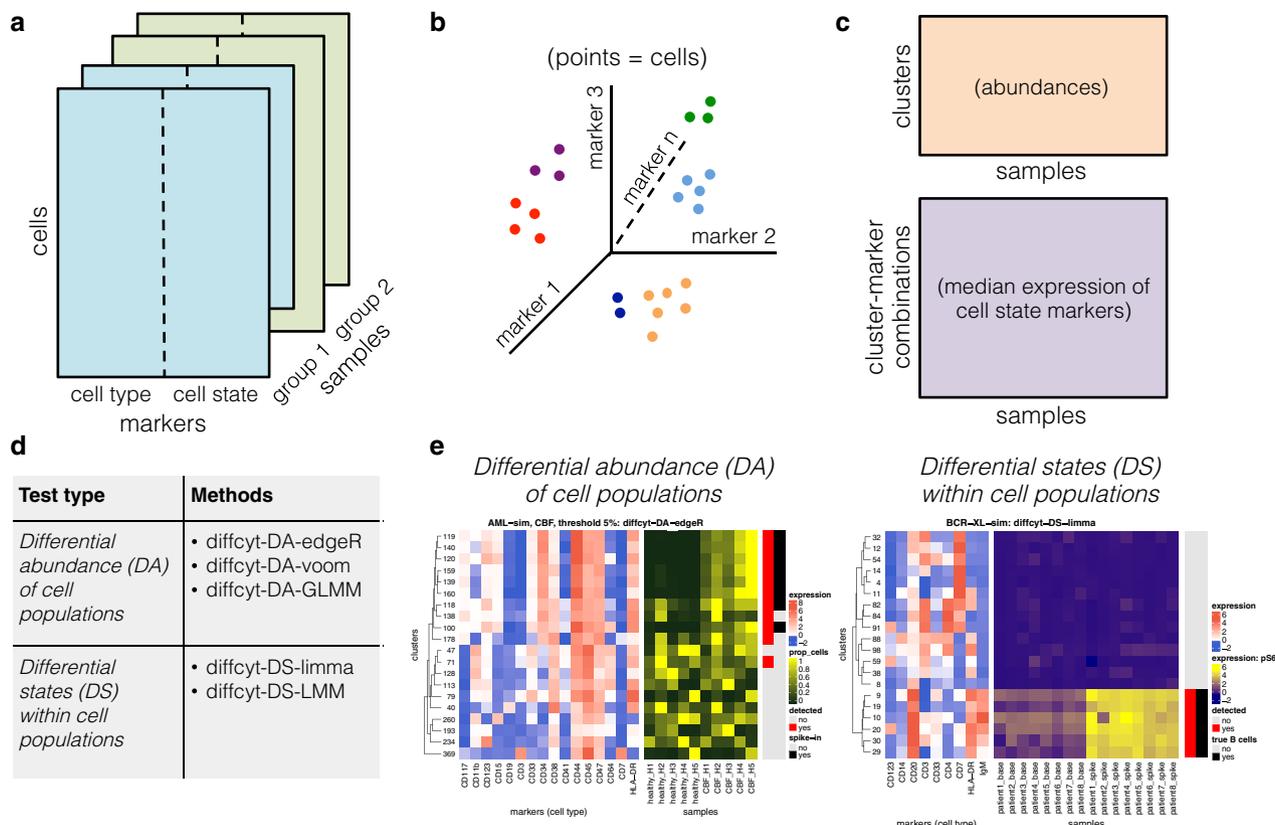


Fig. 1 Schematic overview of *diffcyt* methodology. The *diffcyt* framework applies high-resolution clustering and empirical Bayes moderated tests for differential discovery analyses in high-dimensional cytometry data. **a** Input data are provided as tables of protein marker expression values per cell, one table per sample. Markers may be split into “cell type” and “cell state” categories; in the standard setup, cell type markers are used for clustering. **b** High-resolution clustering summarizes the data into a large number (e.g. 100–400) of clusters representing cell subsets. **c** Features are calculated at the cluster level, including cluster cell counts (abundances), and median expression of cell state markers within clusters. **d** Differential testing methods can be grouped into two types: differential abundance (DA) of cell populations and differential states (DS) within cell populations. Results are returned in the form of adjusted *p*-values, allowing the identification of sets of significant detected clusters (DA tests) or cluster–marker combinations (DS tests). **e** Results are interpreted with the aid of visualizations, such as heatmaps. Example heatmaps show cluster phenotypes (expression profiles) and differential signal of interest (relative cluster abundances or expression of signaling marker pS6, by sample), with annotation for detected significant clusters or cluster–marker combinations (red) and true differential clusters or cluster–marker combinations (black). A detailed description of the *diffcyt* methodology is provided in Methods

cutoffs, and reduced sensitivity at the highest spike-in threshold (5%). *Citrus* detects only a subset of the spiked-in cells, and *cydar* cannot reliably distinguish these rare populations. Figure 2b displays *p*-value distributions from an accompanying null simulation, where no true spike-in signal was included; the *p*-value distributions for the *diffcyt* methods are approximately uniform, indicating good error control and model fit (additional replicates are included in Supplementary Fig. 2). Figure 2c illustrates the expression profiles (phenotypes) and relative abundances by sample for the detected and true differential clusters (additional heatmaps are included in Supplementary Fig. 3). Figure 2d demonstrates the effect of varying the number of clusters across a broad range (between 9 and 1600). Performance is reduced when there are too few clusters (due to merging of populations) or too many clusters (due to low power). The number of clusters is the main parameter choice in the *diffcyt* methods; an optimum is achieved around 400 clusters for this dataset (the remaining thresholds and condition are shown in Supplementary Fig. 4).

Additional results provide further details on overall performance and robustness of the *diffcyt* DA methods. The top detected clusters represent high-precision subsets of the spiked-in population, confirming that the high-resolution clustering strategy has worked as intended (Supplementary Fig. 5). Filtering clusters with low cell counts (using default parameters) did not

remove any clusters from this dataset. An alternative implementation of the *diffcyt-DA-voom* method (using random effects for paired data) gives similar overall performance (Supplementary Fig. 6). Using *FlowSOM* meta-clustering to generate 40 merged clusters instead of testing at high resolution worsens both error control and sensitivity (Supplementary Fig. 7). The influence of random seeds used for the clustering and data generation procedures is greatest at the 0.1% threshold, as expected (Supplementary Figs. 8 and 9). Similarly, additional simulations containing less distinct populations of interest (see Supplementary Note 1) show that reducing signal strength has a strong negative influence on performance at the 0.1% threshold (Supplementary Fig. 10). Using smaller sample sizes (2 vs. 2) affects performance noticeably at the lower thresholds (Supplementary Fig. 11). Finally, runtimes are fastest for methods *diffcyt-DA-edgeR* and *diffcyt-DA-voom* (Supplementary Fig. 12).

Improved performance for DS tests. The second dataset, *BCR-XL-sim*, evaluates performance for detecting DS within cell populations (Figure 3). This dataset contains a spiked-in population of B cells stimulated with B cell receptor/Fc receptor cross-linker (BCR-XL), in a comparison of 8 vs. 8 paired samples of healthy peripheral blood mononuclear cells (original data sourced from ref. 23). The stimulated B cells have elevated expression of

Table 1 Overview of existing methods and limitations

Method	Short description	Limitations	Ref.
<i>Citrus</i>	Uses hierarchical clustering and regularized regression or classification models to select predictive features, such as cluster abundances or median expression of functional markers, that are associated with an outcome of interest	<ul style="list-style-type: none"> Detected features cannot be ranked by importance Lasso-regularized models cannot easily detect multiple correlated features Rare cell populations cannot easily be detected, due to minimum cluster size requirement and computational limitations Response variable is the clinical outcome variable, which makes it difficult to account for complex experimental designs (including batch effects, paired designs, and continuous covariates) 	9
<i>CellCnn</i>	Applies convolutional neural networks in a representation learning framework to detect rare cell populations associated with an outcome of interest; designed specifically for detecting rare cell populations	<ul style="list-style-type: none"> Ranking of detected cells cannot be interpreted in terms of statistical significance Interpretation of detected populations (referred to as filters) can be difficult, since they may be composed of multiple distinct cell populations Response variable is the clinical outcome variable, which makes it difficult to account for complex experimental designs (including batch effects, paired designs, and continuous covariates) All protein markers are treated identically; there is no conceptual split between cell type and cell state (or functional) markers 	10
<i>cydar</i>	Assigns cells to overlapping hyperspheres in the high-dimensional space; tests for differential abundance between hyperspheres using moderated tests from <i>edgeR</i> ^{15,16} , while controlling the spatial false discovery rate among overlapping hyperspheres	<ul style="list-style-type: none"> Rare cell populations cannot easily be detected, due to their relatively small volume in the high-dimensional space All protein markers are treated identically; there is no conceptual split between cell type and cell state (or functional) markers 	11
classic regression-based approach	Automated clustering using <i>FlowSOM</i> ¹⁴ , followed by manual merging and annotation to define cell populations; differential testing of features such as population abundances or median expression of functional markers using generalized linear mixed models, linear mixed models, or linear models	<ul style="list-style-type: none"> Manual merging and annotation step requires expert biological knowledge, and can be time-consuming and subjective When testing large numbers of clusters, e.g. to detect rare cell populations: loss of statistical power due to multiple testing penalty; no sharing of information across clusters 	12

Overview of recently developed methods for performing differential analyses in high-dimensional cytometry data. For each method, a short description of the methodology and a summary of limitations are provided

several signaling state markers, in particular phosphorylated ribosomal protein S6 (pS6); methods are evaluated by their ability to detect differential expression of pS6 within the population of B cells. Figure 3a summarizes performance for the *diffcyt* DS methods and the existing methods. The *diffcyt* methods give the best performance, with *diffcyt-DS-limma* having better error control. *Citrus* and *CellCnn* detect differential expression of pS6 for only a subset of the spiked-in cells, and *cydar* gives poor performance (likely due to ambiguity in assigning cells to overlapping hyperspheres in the high-dimensional space in order to calculate performance metrics). Figure 3b displays *p*-value distributions from a null simulation; *p*-values are approximately uniform across replicates, as previously (additional replicates are included in Supplementary Fig. 13). Figure 3c displays expression profiles of detected and true differential clusters, along with expression by sample of the signaling marker pS6 (additional heatmaps are included in Supplementary Fig. 14). Figure 3d demonstrates the effect of varying the number of clusters. Performance is reduced when there are too few or too many clusters; for this dataset, an optimum is observed across a broad range, including 100 clusters.

As previously, the top detected clusters represent high-precision subsets of the population of interest (Supplementary Fig. 15). Filtering with default parameters did not remove any clusters. To judge the benefit of splitting markers into cell type

and cell state categories, we re-ran the analyses treating all markers as cell type (i.e. used for clustering), and using methods to test for DA instead of DS. This gave similar performance, but makes interpretation more difficult: since the methods test for DA of clusters defined using all markers in this case, the detected differential clusters may mix elements from canonical cell type and cell state phenotypes (Supplementary Fig. 16). Alternative implementations of *diffcyt-DS-limma* (using random effects for paired data) and *diffcyt-DS-LMM* (using fixed effects for paired data) give similar performance overall (Supplementary Fig. 17). For this dataset, using *FlowSOM* meta-clustering to merge clusters does not reduce performance (Supplementary Fig. 18). Varying random seeds for the clustering and data generation procedures does not significantly affect performance (Supplementary Figs. 19 and 20). Additional simulations containing less distinct populations of interest (see Supplementary Note 1) show deteriorating performance when the signal is reduced by 75% (Supplementary Fig. 21). Using smaller sample sizes (4 vs. 4 and 2 vs. 2) worsens error control, especially for *diffcyt-DS-LMM* (Supplementary Fig. 22). Runtimes are fastest for *diffcyt-DS-limma* (Supplementary Fig. 23).

Successful recovery of known signals in experimental data. In order to demonstrate our methods on experimental data, we re-analyzed a dataset from a recent study using mass cytometry to

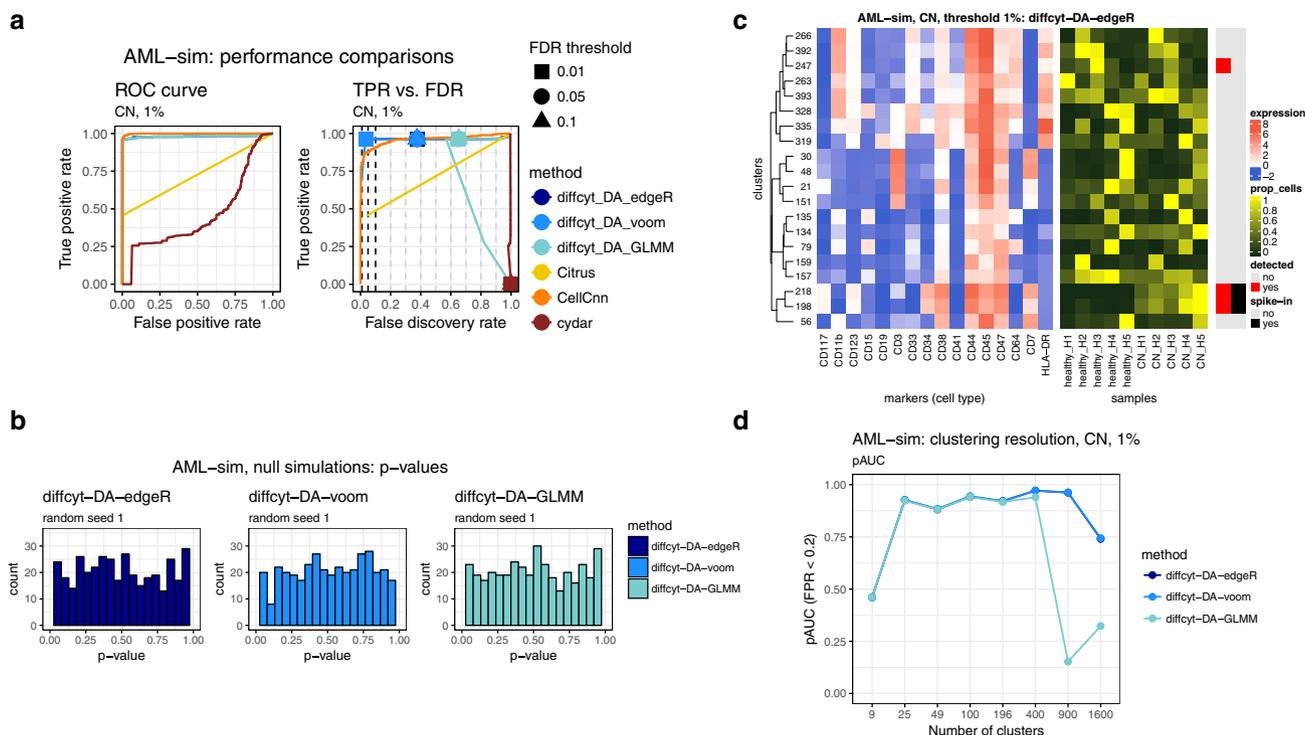


Fig. 2 Benchmarking results for dataset *AML-sim*. **a** Performance metrics for dataset *AML-sim*, testing for differential abundance (DA) of cell populations. Panels show (i) receiver operating characteristic (ROC) curves and (ii) true positive rate (TPR) vs. false discovery rate (FDR) (also indicating observed TPR and FDR at FDR cutoffs 1%, 5%, and 10%). Representative results for one condition (CN vs. healthy) and abundance threshold (1%) are shown (complete results for this dataset are included in Supplementary Fig. 1). **b** Results for additional null simulations, where no true spike-in signal was included; *p*-value distributions are approximately uniform (additional replicates are included in Supplementary Fig. 2). **c** Heatmap displaying phenotypes (expression profiles) of detected and true differentially abundant clusters, along with the signal of interest (relative cluster abundances, by sample), for method *diffcyt-DA-edgeR*. Expression values represent median arcsinh-transformed expression per cluster across all samples (left panel). Rows (clusters) are grouped by hierarchical clustering with Euclidean distance and average linkage; the heatmap shows the top 20 most highly significant clusters. Vertical annotation indicates detected significant clusters at 10% FDR (red) and clusters containing >50% true spiked-in cells (black). (Additional heatmaps are included in Supplementary Fig. 3) **d** Results for varying clustering resolution (between 9 and 1600 clusters), showing partial area under ROC curves (pAUC) for false positive rates (FPR) < 0.2 (additional figures are included in Supplementary Fig. 4). Performance metric plots generated using *iCOBRA*⁴²; heatmaps generated using *ComplexHeatmap*⁴³

characterize immune cell subsets in peripheral blood from melanoma patients treated with anti-PD-1 immunotherapy²⁴ (*Anti-PD-1* dataset; Figure 4). In this study, differential signals were detected for a number of cell populations, both in response to treatment and in baseline comparisons before treatment, between groups of patients classified as responders and non-responders to treatment. One key result was the identification of a small sub-population of monocytes, with frequency in baseline samples (prior to treatment) strongly associated with responder status. The relatively rare frequency made this population difficult to detect; in addition, the dataset contained a strong batch effect due to sample acquisition on two different days²⁴. Using method *diffcyt-DA-edgeR* to perform a differential comparison between baseline samples from the responder and non-responder patients (and taking into account the batch effect), we correctly identified three significant differentially abundant clusters (at an FDR cutoff of 10%) with phenotypes that closely matched the subpopulation of monocytes detected in the original study (CD14⁺ CD33⁺ HLA-DR^{hi} ICAM-1⁺ CD64⁺ CD141⁺ CD86⁺ CD11c⁺ CD38⁺ PD-L1⁺ CD11b⁺ monocytes) (clusters 317, 358, and 380; Figure 4a). One additional cluster with an unknown phenotype was also detected (cluster 308). The total abundance (combined cell counts) of the three matching clusters showed a clear differential signal between the two groups (Figure 4b). However, these results were sensitive to the choice of random seed for the clustering; in five additional runs using different random seeds, we detected

between 0 and 4 significant differentially abundant clusters (at 10% FDR) per run; clusters matching the expected phenotype were detected in four out of the five runs (Supplementary Fig. 24).

For a second evaluation on experimental data, we re-analyzed the original (unmodified) data from the BCR-XL stimulation condition in ref. 23 (*BCR-XL* dataset; Figure 5). This dataset contains strong differential signals for several signaling state markers in several cell populations, as previously described^{12,23}. Using method *diffcyt-DS-limma*, we reproduced several of the major known signals, including strong differential expression of: pS6, pPlcg2, pErk, and pAkt (elevated), and pNFkB (reduced, in BCR-XL stimulated condition) in B cells (identified by expression of CD20); pBtk and pNFkB in CD4⁺ T cells (identified by expression of CD3 and CD4); and pBtk, pNFkB, and pSlp76 in natural killer (NK) cells (identified by expression of CD7). Here, phenotypes can be identified either by marker expression profiles (Figure 5) or, alternatively, using reference population labels available for this dataset (Supplementary Fig. 25).

Discussion

We have presented a new computational framework for performing flexible differential discovery analyses in high-dimensional cytometry data. Our methods are designed for two related but distinct discovery tasks: detecting differentially abundant cell populations, including rare populations; and

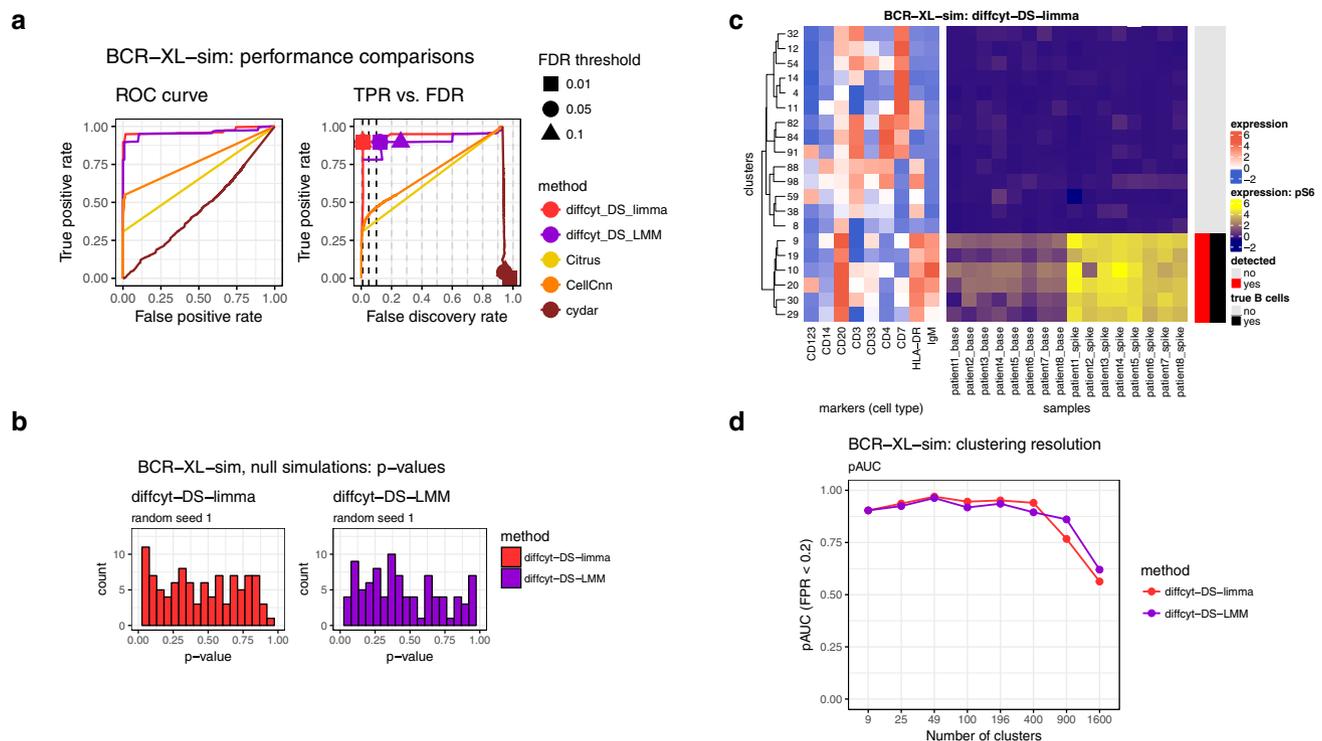


Fig. 3 Benchmarking results for dataset *BCR-XL-sim*. **a** Performance metrics for dataset *BCR-XL-sim*, testing for differential states (DS) within cell populations. Panels show (i) receiver operating characteristic (ROC) curves and (ii) true positive rate (TPR) vs. false discovery rate (FDR) (also indicating observed TPR and FDR at FDR cutoffs 1%, 5%, and 10%). **b** Results for additional null simulations, where no true spike-in signal was included; *p*-value distributions are approximately uniform (additional replicates are included in Supplementary Fig. 13). **c** Heatmap displaying phenotypes (expression profiles) of detected and true differential clusters, along with the signal of interest (expression of signaling marker pS6, by sample), for method *diffcyt-DS-limma*. Expression values represent median arcsinh-transformed expression per cluster across all samples (left panel) or by individual samples (right panel). Rows (clusters) are grouped by hierarchical clustering with Euclidean distance and average linkage; the heatmap shows the top 20 most highly significant clusters. Vertical annotation indicates detected significant cluster-marker combinations at 10% FDR (red) and clusters containing >50% true spiked-in cells (black). (Additional heatmaps are included in Supplementary Fig. 14.) **d** Results for varying clustering resolution (between 9 and 1600 clusters); showing partial area under ROC curves (pAUC) for false positive rates (FPR) < 0.2. Performance metric plots generated using *iCOBRA*⁴²; heatmaps generated using *ComplexHeatmap*⁴³

detecting differential expression of functional or other cell state markers within cell populations. Compared to existing approaches, our methods provide improved detection performance on semi-simulated benchmark datasets, along with fast runtimes. We have also successfully recovered known differential signals in re-analyses of two published experimental datasets, including DA of a highly specific rare population. Our methods can account for complex experimental designs, including batch effects, paired designs, and continuous covariates. In addition, the set of protein markers may be split into cell type and cell state markers, facilitating biological interpretability. Visualizations such as heatmaps can be used to interpret the high-resolution clustering results (for example, to judge whether groups of clusters form larger populations, and to identify the phenotype of detected clusters). Methods *diffcyt-DA-edgeR* (for DA tests) and *diffcyt-DS-limma* (for DS tests) achieved the best performance and fastest runtimes overall (Figure 2 and 3); we recommend these as the default choices.

One limitation of our framework is that groups of similar clusters cannot be automatically merged into larger cell populations with a consistent phenotype. For example, the clear group of detected clusters in Figure 3c would ideally be merged into a single population representing B cells. However, this is a difficult computational problem, since the optimal resolution depends on the biological setting, and any automatic merging must avoid merging rare cell populations into larger ones. Our high-resolution clustering approach instead provides a tractable

“middle ground” between discrete clustering and a continuum of cell populations; we return results directly at the level of high-resolution clusters, and let the user interpret them via visualizations. A related issue concerns the identification of cell population phenotypes: our approach relies on visualizations and manual annotation of populations, which necessarily involves some subjectivity. Recently, several new methods have been published for automated labeling of cell populations²⁵, identification of simplified gating strategies to describe cell populations of interest^{26,27}, or to compare cluster phenotypes²⁸. These methods could be integrated within our framework to interpret detected differential clusters in a more automated manner. Similarly, clustering algorithms that generate biologically interpretable clusters could be used to improve interpretability²⁹.

A further limitation relates to batch effects: in datasets with strong batch effects, the high-resolution clustering may separate across batches, making it more difficult to distinguish the signal of interest. Aligning cell populations across batches is an active area of research in single-cell analysis (e.g. refs. 30–33); ideally, these methods will be integrated with frameworks for downstream differential analyses. Another issue concerns our strategy of summarizing cell state marker signals into median values. This strategy has advantages of simplicity, ease of interpretation, and fast runtimes. However, some information is necessarily lost, especially for markers with multi-modal distributions; good frameworks for flexible comparisons of full distributions are currently lacking. Additionally, splitting markers into groups

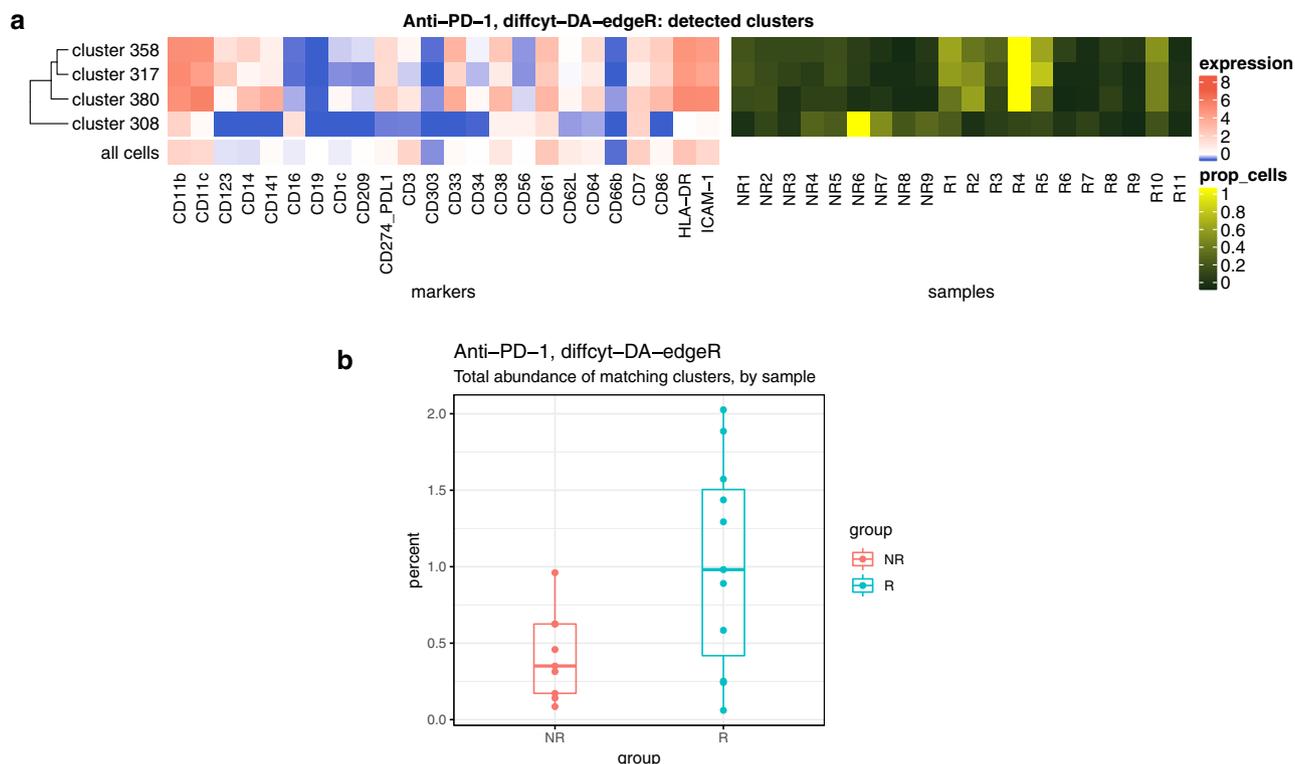


Fig. 4 Results for experimental dataset *Anti-PD-1*. Results for re-analysis of experimental dataset *Anti-PD-1* using method *diffcyt-DA-edgeR*; testing for differential abundance (DA) of cell populations between baseline samples from responder and non-responder groups of patients. **a** Heatmap shows phenotype (median arcsinh-transformed marker expression profiles) of significant detected clusters at 10% false discovery rate (FDR), compared to all cells (left panel); and relative cluster abundances (proportion of cells per cluster, by sample) (right panel) for the detected clusters. Heatmap rows (clusters) are grouped by hierarchical clustering with Euclidean distance and average linkage. **b** Boxplot shows total abundance (combined number of cells) for the clusters matching the phenotype of interest (clusters 317, 358, and 380), by sample and group. Runtime was 32.0 s, on a 2014 MacBook Air laptop, 1.7 GHz processor, 8 GB memory, using a single processor core. NR non-responders, R responders

representing cell type and cell state may be seen as a disadvantage in applications where this distinction is not clear. However, this step is optional: it is possible to run our methods using all markers for clustering (i.e. treating all markers as cell type) and testing for DA (Supplementary Fig. 16). For well-characterized immune populations, standard cell type markers may be found in the literature (e.g. ref. ³⁴) or by consulting the websites of commercial antibody suppliers (e.g. BioLegend, Miltenyi Biotec, or Bio-Rad). Methods are also available to automatically group markers^{12,22,35}, although these should be used with care to ensure that cell population definitions are biologically plausible. For markers with subtle shifts (e.g. cytokines), assigning these as cell state markers and applying DS tests may fail to detect the differential signal; in this case, cluster labels may be exported to facilitate alternative analysis strategies (e.g. visualizations using *CytoRSuite*³⁶, *iSEE*³⁷, *OpenCyto*³⁸, or commercial software such as *FlowJo*).

The main user parameter in our methods is the number of clusters. The optimal value depends on several factors, including the size of the dataset (number of cells and samples), the expected relative abundances of cell populations of interest, and the number of markers used to define cell populations. The number of clusters determines the number of statistical tests, and affects power through the multiple testing penalty and the counts per cluster. We recommend higher numbers of clusters when rare cell populations are of interest (for example, we used 400 clusters for the *AML-sim* dataset, and 100 clusters for the *BCR-XL-sim* dataset). Ultimately, this is a subjective choice for the user, which may also be explored interactively: e.g. by trying several different resolutions, and judging the interpretability of the results using

visualizations or by calculating cluster separation metrics (e.g. average silhouette width). However, in our evaluations, good results were obtained over a range of resolutions (Figure 2d and 3d). Most computational methods include one or more parameters that can be adjusted by the user; in our view, one of the advantages of our approach is that the number of clusters is an intuitive parameter, with values that can be easily interpreted.

In general, we note that our methods are designed for “discovery” analyses: all results should be explored and interpreted using visualizations, and any generated hypotheses must ultimately be validated with targeted confirmatory experiments. Our methods are implemented in the open-source R package *diffcyt*, available from Bioconductor (<http://bioconductor.org/packages/diffcyt>). The package includes comprehensive documentation and code examples, including an extended workflow vignette. Code to reproduce all analyses and figures from our benchmarking evaluations is available from GitHub (<https://github.com/lmweber/diffcyt-evaluations>), and data files from the benchmarking datasets are available from FlowRepository³⁹ (<http://flowrepository.org/id/FR-FCM-ZYL8>), allowing other researchers to extend and build on our analyses.

Methods

Description of *diffcyt* methodology. The following sections provide a detailed description of the *diffcyt* methodology (see Figure 1 for a schematic overview).

Preprocessing. *Data preparation:* Input data is formatted into a Bioconductor *SummarizedExperiment* object containing a single matrix of protein expression values, with one row per cell, and one column per protein marker. Row meta-data contains sample IDs and group IDs, and column meta-data contains protein marker information. The *SummarizedExperiment* format enables easy subsetting

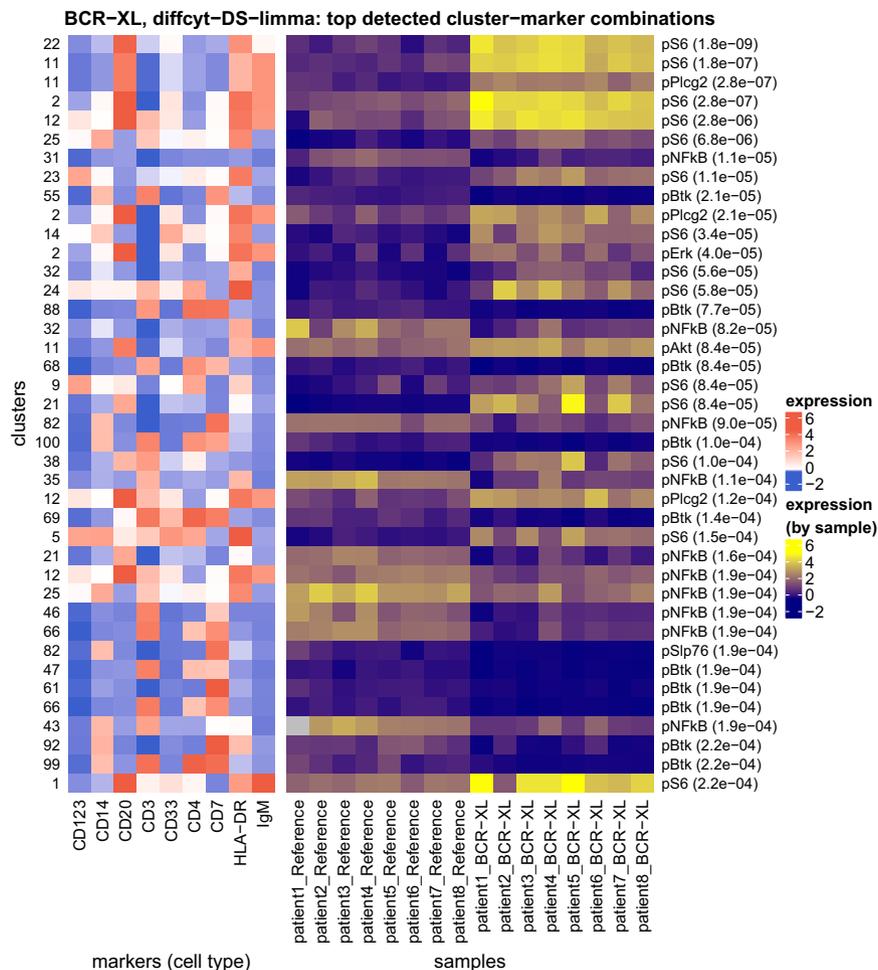


Fig. 5 Results for experimental dataset BCR-XL. Results for re-analysis of experimental dataset *BCR-XL* using method *diffcyt-DS-limma*; testing for differential states (DS) within cell populations. Heatmap shows phenotypes (median arcsinh-transformed expression profiles for cell type markers) for the top 40 most highly significant detected cluster-marker combinations (left panel), and expression by sample for the cell state marker (signaling marker) in each detected cluster-marker combination (right panel). Rows (cluster-marker combinations) are ordered by decreasing adjusted *p*-values. Cell state marker names and adjusted *p*-values are displayed in right-hand-side row headings. Color scale for expression of cell type markers is normalized to 1st and 99th percentiles across all clusters and markers. Only the top 40 most highly significant detected cluster-marker combinations (out of 1400 total) are shown, for easier visibility. Runtime was 15.0 s, on a 2014 MacBook Air laptop, 1.7 GHz processor, 8 GB memory, using a single processor core

of data and meta-data, as well as simplified interaction with other Bioconductor packages.

Marker information: “cell type” and “cell state” markers: The set of protein markers may be split into sets of “cell type” and “cell state” markers. This split enables the methodology to take advantage of existing biological knowledge, and facilitates interpretability. By default, cell type markers are used to define clusters representing cell populations (which are tested for DA), and median cell state marker signals are used to test for DS (e.g. signaling or other functional states) within populations. This allows the user to interpret the results in terms of cell populations defined by known cell type markers.

The grouping into cell type and cell state markers must be specified by the user, and is stored in the column meta-data of the *SummarizedExperiment* object. This grouping is an important design choice, which may be made based on prior biological knowledge or using data-driven methods. For an example of a data-driven method of marker ranking and selection, see refs. 22,12.

Subsampling: Optionally, random subsampling can be used to select an equal number of cells from each sample. This can be useful when there are large differences in total numbers of cells per sample, since it ensures that samples with relatively large numbers of cells do not dominate the clustering. However, some information will necessarily be lost. Subsampling should generally not be used when rare cell populations are of interest, due to the significant loss of information if cells from the rare population are discarded.

Transformation: Expression values are transformed using an inverse hyperbolic sine (*arcsinh*) transform with adjustable *cofactor* parameter. Raw expression values (fluorescence intensities for flow cytometry, or randomized ion counts for mass cytometry) follow an approximately log-normal distribution; the *arcsinh* transform brings this closer to a normal distribution (or mixture of normal distributions), which improves clustering performance and allows positive and negative

populations to be distinguished more clearly. The *arcsinh* transform behaves similarly to a log transform at high values, but is approximately linear near zero; so unlike the log, it can handle zeros or small negative values. The *cofactor* controls the width of the linear region. (Zero values and small negatives occur in mass cytometry data when no ions are detected in a given channel: negatives are due to background subtraction and randomization of integer count values, which are performed by default by the instrument software.) Standard values for the *cofactor* are 5 for mass cytometry, and 150 for flow cytometry (see ref. 2, Supplementary Fig. S2).

Integration with CATALYST package: Alternatively, a pre-prepared *daFrame* object from the *CATALYST* R/Bioconductor package⁴⁰ can be used as the input for the *diffcyt* methods. The *CATALYST* package contains extensive functions for preprocessing, exploratory analysis, and visualization of mass cytometry data. If this option is used, preprocessing (and clustering) are done using *CATALYST*. This is particularly useful when *CATALYST* has already been used for exploratory analyses and visualizations; the *diffcyt* package can then be used to calculate differential tests. For more details, see the *diffcyt* and *CATALYST* Bioconductor package vignettes.

Clustering. The clustering step is the core of the *diffcyt* methodology. We use high-resolution clustering to group cells into a large number of small clusters representing cell populations or subsets, which can then be further analyzed by differential testing. High-resolution clustering (or over-clustering) helps ensure that small or rare cell populations are adequately separated from larger populations.

By default, we use the *FlowSOM* clustering algorithm¹⁴ (available from Bioconductor) to generate the clusters, since we previously showed that *FlowSOM* gives very good clustering performance for high-dimensional cytometry data, for both major and rare cell populations, and is extremely fast⁷. However, we run

FlowSOM without the final meta-clustering step, to help ensure that small or rare populations are not merged into larger populations, which is crucial for detecting DA of extremely rare populations.

If markers have been split into sets of cell type and cell state markers, then (by default) the clustering is performed using cell type markers only.

Data features. After clustering, we calculate features summarizing the data at the cluster level: cluster cell counts or abundances (number of cells per cluster–sample combination), and median transformed marker expression values (per cluster–sample combination). The feature values are formatted as new *SummarizedExperiment* objects, where rows represent clusters or cluster–marker combinations, and columns represent samples. These feature values are then used as inputs for the differential testing.

Design matrices and model formulas. The models to be fitted are specified with a design matrix or model formula, depending on the differential testing method used. Design matrices consist of one row per sample, and columns containing predictor variables, including the outcome of interest (e.g. columns of indicator variables for group IDs, such as diseased and healthy) and any other covariates. Flexible experimental designs are possible: block IDs (e.g. patient IDs in a paired design), batch effects, and continuous covariates can be included in the design matrix; each of these terms will be included as fixed effects in the models. Alternatively, model formulas also provide the option to include block IDs as random intercept terms (instead of fixed effects). When testing for DA, model formulas can also be used to include random intercept terms for each sample (known as “observation-level random effects” or OLRs; see ref. 12), to account for overdispersion typically seen in high-dimensional cytometry data.

Contrasts. The comparison of interest for the differential tests is specified with a contrast matrix. The contrast matrix consists of one row per model coefficient (corresponding to columns from the design matrix), and a column specifying the comparison of interest (i.e. the combination of model coefficients that is assumed to equal zero under the null hypothesis). This system of combining a design matrix (or model formula) with an appropriate contrast matrix provides users with powerful options to investigate a wide range of possible hypotheses within flexible experimental design settings.

Tests for DA of cell populations. *diffcyt-DA-edgeR*: The *diffcyt-DA-edgeR* method calculates tests for DA of clusters using methodology from the *edgeR* package^{15,16}. This method uses *edgeR* to fit models and calculate moderated tests at the cluster level. The moderated tests improve power by sharing information on variability (i.e. variance across samples for a single cluster) across clusters. Note that by default, we use the option *trend.method = “none”* to estimate common dispersions (see *edgeR* User’s Guide, available from Bioconductor).

The input to the tests is a table of cluster cell counts. The experimental design is specified using a design matrix, which enables flexible experimental designs. The comparison of interest is specified using a contrast matrix. A filtering step removes clusters with very low cell counts across samples to improve power. Normalization for the total number of cells per sample (library sizes) is automatically performed by the *edgeR* functions. Optionally, normalization factors for composition effects can be calculated using the “trimmed mean of *M*-values” (TMM) method from the *edgeR* package⁴¹.

Differential test results are returned in the form of raw *p*-values and adjusted *p*-values (FDR) from the moderated tests, which can be used to rank the clusters by their evidence for DA. The results are stored in a new *SummarizedExperiment* object.

diffcyt-DA-voom: The *diffcyt-DA-voom* method calculates tests for DA of clusters using methodology from the *limma* package¹⁷ and *voom* method¹⁸. This method uses *limma* to fit models and calculate moderated tests at the cluster level. The moderated tests improve power by sharing information on variability across clusters. Since count data (such as cluster cell counts) are often heteroscedastic, we use *voom* to transform the raw cluster cell counts and estimate observation-level precision weights in order to stabilize the mean–variance relationship.

The input to the tests is a table of cluster cell counts. The experimental design is specified using a design matrix, which enables flexible experimental designs. For paired designs, either fixed effects or random effects can be used; fixed effects are simpler, but random effects may improve power in datasets with unbalanced designs or very large numbers of samples. Random effects make use of the *limma duplicateCorrelation* methodology (note that this methodology does not allow multiple measures per sample; in this case, fixed effects should be used instead). The comparison of interest is specified using a contrast matrix. A filtering step removes clusters with very low cell counts across samples to improve power. Normalization for the total number of cells per sample (library sizes) is automatically performed by the *limma* and *voom* functions. Optionally, normalization factors for composition effects can be calculated using the TMM method from the *edgeR* package⁴¹.

Differential test results are returned in the form of raw *p*-values and adjusted *p*-values (FDR) from the moderated tests, which can be used to rank the clusters by

their evidence for DA. The results are stored in a new *SummarizedExperiment* object.

diffcyt-DA-GLMM: The *diffcyt-DA-GLMM* method calculates tests for DA of clusters using the generalized linear mixed models (GLMM) methodology originally implemented by ref. 12. This method fits GLMMs for each cluster, and calculates differential tests separately for each cluster (i.e. one model per cluster). The response variables in the models are the cluster cell counts, which are assumed to follow a binomial distribution. Note that the original methodology from ref. 12 has been modified here to make use of high-resolution clustering to enable rare cell populations to be investigated more easily. In addition, we do not attempt to manually merge clusters into canonical cell populations; results are instead reported directly at the high-resolution cluster level.

The input to the tests is a table of cluster cell counts. The experimental design is specified using a model formula, which enables flexible experimental designs. Blocking variables (e.g. for paired designs) can be included as either random intercept terms or fixed effect terms. For paired designs, we recommend using random intercept terms to improve statistical power (see ref. 12). Batch effects and continuous covariates are included as fixed effects. In addition, we include random intercept terms for each sample to account for overdispersion typically seen in high-dimensional cytometry count data. The sample-level random intercept terms are known as “observation-level random effects” (OLREs; see ref. 12). The comparison of interest is specified using a contrast matrix. A filtering step removes clusters with very low cell counts across samples to improve power. Optionally, normalization factors for composition effects can be calculated using the TMM method from the *edgeR* package⁴¹.

Differential test results are returned in the form of raw *p*-values and adjusted *p*-values (FDR), which can be used to rank the clusters by their evidence for DA. The results are stored in a new *SummarizedExperiment* object.

Tests for DS within cell populations. *diffcyt-DS-limma*: The *diffcyt-DS-limma* method calculates tests for DS within clusters using methodology from the *limma* package¹⁷. Clusters are defined using cell type markers, and cell states are defined using median transformed expression of cell state markers within clusters. This method uses *limma* to fit models and calculate moderated tests at the cluster level. The moderated tests improve power by sharing information on variability across clusters. Note that by default, we use the option *trend = TRUE* in the *limma eBayes* fitting function in order to stabilize the mean–variance relationship.

The input to the tests is a set of tables of median expression of each marker for each cluster–sample combination. The experimental design is specified using a design matrix, which enables flexible experimental designs. For paired designs, either fixed effects or random effects can be used; fixed effects are simpler, but random effects may improve power in datasets with unbalanced designs or very large numbers of samples. Random effects make use of the *limma duplicateCorrelation* methodology (note that this methodology does not allow multiple measures per sample; in this case, fixed effects should be used instead). The comparison of interest is specified using a contrast matrix. A filtering step removes clusters with very low cell counts across samples to improve power. If cluster cell counts are provided, these can be used to calculate precision weights (across all samples and clusters), allowing the *limma* model fitting functions to account for uncertainty due to the total number of cells per sample (library size normalization) and total number of cells per cluster.

Differential test results are returned in the form of raw *p*-values and adjusted *p*-values (FDR) from the moderated tests for each cluster–marker combination (for cell state markers). These can be used to rank the cluster–marker combinations by their evidence for DS. The results are stored in a new *SummarizedExperiment* object.

diffcyt-DS-LMM: The *diffcyt-DS-LMM* method calculates tests for DS within clusters using the linear mixed models (LMM) and linear models (LM) methodology originally implemented by ref. 12. Clusters are defined using cell type markers, and cell states are defined using median transformed expression of cell state markers within clusters. This method fits LMMs for each cluster–marker combination (for cell state markers), and calculates differential tests separately for each cluster–marker combination (i.e. one model per cluster–marker combination). The response variable in each model is the median arcsinh-transformed marker expression of the cell state marker, which is assumed to follow a normal distribution. Note that the original methodology from ref. 12 has been modified here to make use of high-resolution clustering to enable rare cell populations to be investigated more easily. In addition, we do not attempt to manually merge clusters into canonical cell populations; results are instead reported directly at the high-resolution cluster level.

The input is a set of tables of median expression of each marker for each cluster–sample combination. The experimental design is specified using a model formula, which enables flexible experimental designs. Blocking variables (e.g. for paired designs) can be included as either random intercept terms or fixed effect terms. For paired designs, we recommend using random intercept terms to improve statistical power (see ref. 12). Batch effects and continuous covariates are included as fixed effects. If no random intercept terms are included in the model formula, model fitting is performed using an LM instead of an LMM. The comparison of interest is specified using a contrast matrix. A filtering step removes clusters with very low cell counts across samples to improve power. Within each

model, sample-level weights can be included for the number of cells per sample; these weights represent the relative uncertainty in calculating each median value. (Additional uncertainty exists due to variation in the total number of cells per cluster; however, it is not possible to account for this, since separate models are used for each cluster–marker combination.)

Differential test results are returned in the form of raw p -values and adjusted p -values (FDR) for each cluster–marker combination (for cell state markers). These can be used to rank the cluster–marker combinations by their evidence for DS. The results are stored in a new *SummarizedExperiment* object.

Interpretation and visualization. The *diffcyt* methods return results in the form of adjusted p -values (FDR) at the level of high-resolution clusters, either for a given cluster (for DA tests) or cluster–marker combination (for DS tests).

Due to the high-resolution clustering strategy, detected differential cell populations may be split into several sub-clusters with similar phenotypes. For biological interpretation, it is often useful to group the high-resolution clusters into larger populations with a consistent phenotype. However, automatically aggregating clusters is a difficult computational task, since the optimal resolution depends on the biological setting. In particular, there is a risk of merging rare cell populations into larger populations. Therefore, we have adopted the approach of returning results directly at the high-resolution cluster level. These results can then be explored and interpreted using visualizations.

Detailed visualizations can be generated using plotting functions from the *CATALYST* R/Bioconductor package⁴⁰, which accepts output objects from *diffcyt*. Key visualizations include heatmaps showing the phenotype (marker expression profiles) of detected clusters together with the sample-level signal of interest (cluster abundance or median expression of cell state markers). Examples are provided in the *diffcyt* and *CATALYST* Bioconductor package vignettes.

Number of clusters. The number of clusters is the main user parameter choice in the *diffcyt* methods. In the default implementation using the *FlowSOM* algorithm for clustering, this can be specified with the two arguments $xdim$ and $ydim$ in the function *generateClusters*. The total number of clusters is then $xdim * ydim$. (This format is required since *FlowSOM* arranges clusters in a two-dimensional self-organizing map grid.)

The default is 100 clusters ($xdim = 10$, $ydim = 10$), which we expect is sufficient for many datasets. In general, we recommend higher numbers of clusters for datasets where rare cell populations are of interest. In our benchmarking evaluations, we used 400 clusters for the *AML-sim* dataset and 100 clusters for the *BCR-XL-sim* dataset. Ultimately, this is a subjective choice for the user, which will depend on the biological setting and questions of interest in a given dataset; strategies to determine an appropriate number may include interactive exploration of visualizations, and (if available) making use of manually gated populations as a reference.

Benchmark datasets. A complete description of the benchmark datasets used to evaluate the methods is provided in Supplementary Note 1 (including Supplementary Figs. 26 and 27).

Comparisons with existing methods. Additional details on the comparisons with existing methods are provided in Supplementary Note 2.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

Data files for all benchmark datasets are available in FCS format from FlowRepository³⁹ (repository ID: FR-FCM-ZYL8) at <http://flowrepository.org/id/FR-FCM-ZYL8>. The benchmark datasets can also be accessed in *SummarizedExperiment* and *flowSet* Bioconductor object formats through the *HDCytoData* Bioconductor package, available at <http://bioconductor.org/packages/HDCytoData>.

Code availability

The methods described in this paper are implemented in the open-source R package *diffcyt*, which is freely available from Bioconductor at <http://bioconductor.org/packages/diffcyt>. The *diffcyt* package includes comprehensive help files for each function, as well as a package vignette demonstrating a complete example workflow. Code scripts to reproduce all performance evaluations and comparisons with existing methods, reproduce all data preparation and simulation steps, and generate all figures are available from GitHub at <https://github.com/lmweber/diffcyt-evaluations>. The results and figures in this paper were generated using *diffcyt* version 1.3.0 (available from GitHub at <https://github.com/lmweber/diffcyt/releases>) and R version 3.5.0.

Received: 20 December 2018 Accepted: 5 April 2019

Published online: 14 May 2019

References

- Saeyns, Y., Van Gassen, S. & Lambrecht, B. N. Computational flow cytometry: helping to make sense of high-dimensional immunology data. *Nat. Rev. Immunol.* **16**, 449–462 (2016).
- Bendall, S. C. et al. Single-cell mass cytometry of differential immune and drug responses across a human hematopoietic continuum. *Science* **332**, 687–696 (2011).
- Shahi, P., Kim, S. C., Haliburton, J. R., Gartner, Z. J. & Abate, A. R. Abseq: ultrahigh-throughput single cell protein profiling with droplet microfluidic barcoding. *Sci. Rep.* **7**, 44447 (2017).
- Stoeckius, M. et al. Simultaneous epitope and transcriptome measurement in single cells. *Nat. Methods* **14**, 865–868 (2017).
- Peterson, V. M. et al. Multiplexed quantification of proteins and transcripts in single cells. *Nat. Biotechnol.* **35**, 936–939 (2017).
- Aghaeepour, N. et al. Critical assessment of automated flow cytometry data analysis techniques. *Nat. Methods* **10**, 228–238 (2013).
- Weber, L. M. & Robinson, M. D. Comparison of clustering methods for high-dimensional single-cell flow and mass cytometry data. *Cytometry A* **89A**, 1084–1096 (2016).
- Aghaeepour, N. et al. A benchmark for evaluation of algorithms for identification of cellular correlates of clinical outcomes. *Cytometry A* **89A**, 16–21 (2016).
- Bruggner, R. V., Bodenmiller, B., Dill, D. L., Tibshirani, R. J. & Nolan, G. P. Automated identification of stratifying signatures in cellular subpopulations. *Proc. Natl Acad. Sci. USA* **111**, E2770–E2777 (2014).
- Arvaniti, E. & Claassen, M. Sensitive detection of rare disease-associated cell subsets via representation learning. *Nat. Commun.* **8**, 1–10 (2017).
- Lun, A. T. L., Richard, A. C. & Marioni, J. C. Testing for differential abundance in mass cytometry data. *Nat. Methods* **14**, 707–709 (2017).
- Nowicka, M. et al. CyTOF workflow: differential discovery in high-throughput high-dimensional cytometry datasets. *F1000Res.* **6**, 748, version 2 (2017).
- Fonseka, C. Y. et al. Mixed-effects association of single cells identifies an expanded effector CD4+ T cell subset in rheumatoid arthritis. *Sci. Transl. Med.* **10**, eaq0305 (2018).
- Van Gassen, S. et al. FlowSOM: using self-organizing maps for visualization and interpretation of cytometry data. *Cytometry A* **87A**, 636–645 (2015).
- Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–140 (2010).
- McCarthy, D. J., Chen, Y. & Smyth, G. K. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res.* **40**, 4288–4297 (2012).
- Ritchie, M. E. et al. limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
- Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol.* **15**, R29 (2014).
- Wagner, A., Regev, A. & Yosef, N. Revealing the vectors of cellular identity with single-cell genomics. *Nat. Biotechnol.* **34**, 1145–1160 (2016).
- Regev, A. et al. and Human Cell Atlas Meeting Participants. The Human Cell Atlas. *eLIFE* **6**, 1–30 (2017).
- Zeng, H. & Sanes, J. R. Neuronal cell-type classification: challenges, opportunities and the path forward. *Nat. Rev. Neurosci.* **18**, 530–546 (2017).
- Levine, J. H. et al. Data-driven phenotypic dissection of AML reveals progenitor-like cells that correlate with prognosis. *Cell* **162**, 184–197 (2015).
- Bodenmiller, B. et al. Multiplexed mass cytometry profiling of cellular states perturbed by small-molecule regulators. *Nat. Biotechnol.* **30**, 858–867 (2012).
- Krieg, C. et al. High-dimensional single-cell analysis predicts response to anti-PD-1 immunotherapy. *Nat. Med.* **24**, 144–153 (2018).
- Abdelal, T. et al. Predicting cell populations in single cell mass cytometry data. *Cytometry A* <https://doi.org/10.1002/cyto.a.23738> (2019).
- Becht, E. et al. Reverse-engineering flow-cytometry gating strategies for phenotypic labelling and high-performance cell sorting. *Bioinformatics* **35**, 301–308 (2018).
- Aghaeepour, N. et al. GateFinder: projection-based gating strategy optimization for flow and mass cytometry. *Bioinformatics* **34**, 4131–4133 (2018).
- Platon, L. et al. A computational approach for phenotypic comparisons of cell populations in high-dimensional cytometry data. *Methods* **132**, 66–75 (2018).
- Commenges, D., Alkassim, C., Gottardo, R., Hejblum, B. & Thiébaud, R. cytometree: a binary tree algorithm for automatic gating in cytometry analysis. *Cytometry A* **93A**, 1132–1140 (2018).
- Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.* **36**, 411–420 (2018).
- Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat. Biotechnol.* **36**, 421–427 (2018).
- Orlova, D. Y. et al. QFMatch: multidimensional flow and mass cytometry samples alignment. *Sci. Rep.* **8**, 1–14 (2018).

33. Li, Y. H. et al. Scalable multi-sample single-cell data analysis by partition-assisted clustering and multiple alignments of networks. *PLoS Comput. Biol.* **13**, 1–37 (2017).
34. Engel, P. et al. CD nomenclature 2015: human leukocyte differentiation antigen workshops as a driving force in immunology. *J. Immunol.* **195**, 4555–4563 (2015).
35. Diggins, K. E., Greenplate, A. R., Leelatian, N., Wogslund, C. E. & Irish, J. M. Characterizing cell subsets using marker enrichment modeling. *Nat. Methods* **14**, 275–278 (2017).
36. Hammill, D. CytoRSuite. R package, version 0.9.9, <https://github.com/DillonHammill/CytoRSuite> (2019).
37. Rue-Albrecht, K., Marini, F., Sonesson, C. & Lun, A. T. L. iSEE: Interactive SummarizedExperiment Explorer. *F1000Res.* **7**, 741 (2018).
38. Finak, G. et al. OpenCyto: an open source infrastructure for scalable, robust, reproducible, and automated, end-to-end flow cytometry data analysis. *PLoS Comput. Biol.* **10**, e1003806 (2014).
39. Spidlen, J., Breuer, K., Rosenberg, C., Kotecha, N. & Brinkman, R. R. FlowRepository: a resource of annotated flow cytometry datasets associated with peer-reviewed publications. *Cytometry A* **81A**, 727–731 (2012).
40. Chevrier, S. et al. Compensation of signal spillover in suspension and imaging mass cytometry. *Cell Syst.* **6**, 612–620 (2018).
41. Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol.* **11**, R25 (2010).
42. Sonesson, C. & Robinson, M. D. iCOBRA: open, reproducible, standardized and live method benchmarking. *Nat. Methods* **13**, 283 (2016).
43. Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in multidimensional genomic data. *Bioinformatics* **32**, 2847–2849 (2016).

Acknowledgements

The authors thank Helena L. Crowell (University of Zurich) for feedback on the implementation of the *diffcyt* R package, and all members of the Robinson Lab at the University of Zurich for feedback on the methodology and benchmarking. L.M.W. was supported by a Forschungskredit (Candoc) grant from the University of Zurich

(FK-17-100). M.D.R. acknowledges support from the University Research Priority Program Evolution in Action at the University of Zurich.

Author contributions

L.M.W. and M.D.R. developed methods, designed analyses, and wrote the manuscript. L.M.W. implemented methods and performed analyses. M.N. developed methods and assisted with interpretation. C.S. assisted with designing analyses and interpretation. All authors read and approved the final manuscript.

Additional information

Supplementary information accompanies this paper at <https://doi.org/10.1038/s42003-019-0415-5>.

Competing interests: The authors declare no competing interests.

Reprints and permission information is available online at <http://npg.nature.com/reprintsandpermissions/>

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2019